

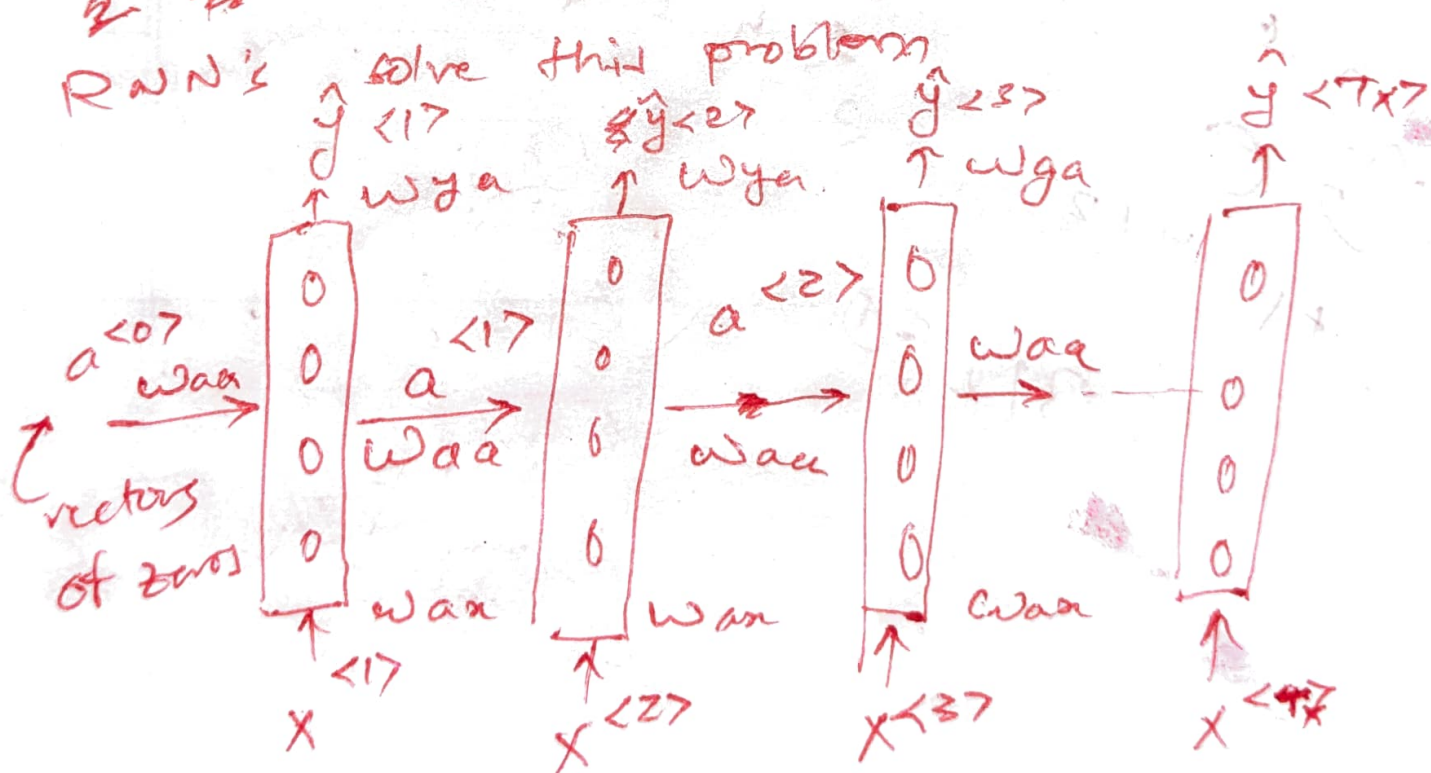
Course 5 (RNN & LSTM)

Why we cannot use standard NN for NLP tasks?

1. Input size in standard NN is fixed, whereas, no. of words in a sentence may vary.
2. Even if we pad the ~~left~~ input which is empty (say standard input size is 100, no. of words are 70 & mark the ~~left~~ ^{empty} 30) with 0, it will require lot of computational power. (Due to one-hot encoding)

2. ~~Q~~

RNN's solve this problem



Wax, Wxy, Waa are shared across all x. i.e. its same for all x.

The output of the previous input is fed ^{along with} current input. i.e. prediction is only based on ^{on} ~~words~~ input words in earlier sequence but not information later in the sequence. This is a drawback of this particular RNN architecture.

Forward ~~prop~~ propagation

$$a^{<0>} = 0, \quad a^{<1>} = g^{(1)}(w_{aa} a^{<0>} + w_{ax} \cdot x^{<1>} + b_a)$$

$$y^{<1>} = g^{(2)}(w_{ya} a^{<1>} + b_y)$$

$g^{(1)} \rightarrow \tanh, \text{ReLU}$

for binary classification \rightarrow sigmoid

for k-way classification \rightarrow softmax

General representation

$$\textcircled{1} \quad a^{<t>} = g(w_{aa} a^{<t-1>} + w_{ax} x^{<t>} + b_a)$$

$$\textcircled{2} \quad y^{<t>} = g(w_{ay} a^{<t>} + b_y)$$

Expression ① can be written as

$$a^{<t>} = g(w_a(a^{<t-1>}, x^{<t>})) + b_a$$

w_a will be composed of w_{aa} & w_{ax}

Say $w_{aa} = (100, 100)$

$w_{ax} = (1000, 10,000)$

Then $w_a = [w_{aa} \mid w_{ax}]$

$\leftarrow 100 \rightarrow \leftarrow 10,000 \rightarrow$

$\in \mathbb{R}^{1 \times 2}$

$= (100, 10,000)$

$(a^{<t-1>}, x^{<t>}) = \left[\frac{a^{<t-1>}}{x^{<t-1>}} \right]$

$\leftarrow 100 \rightarrow$
 $\leftarrow 10,000 \rightarrow$

$\in \mathbb{R}^{2 \times 1}$

$$[W_{aa} \quad W_{an}] \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} = W_{aa} \cdot a^{<t-1>} + W_{ax} \cdot x^{<t>}$$

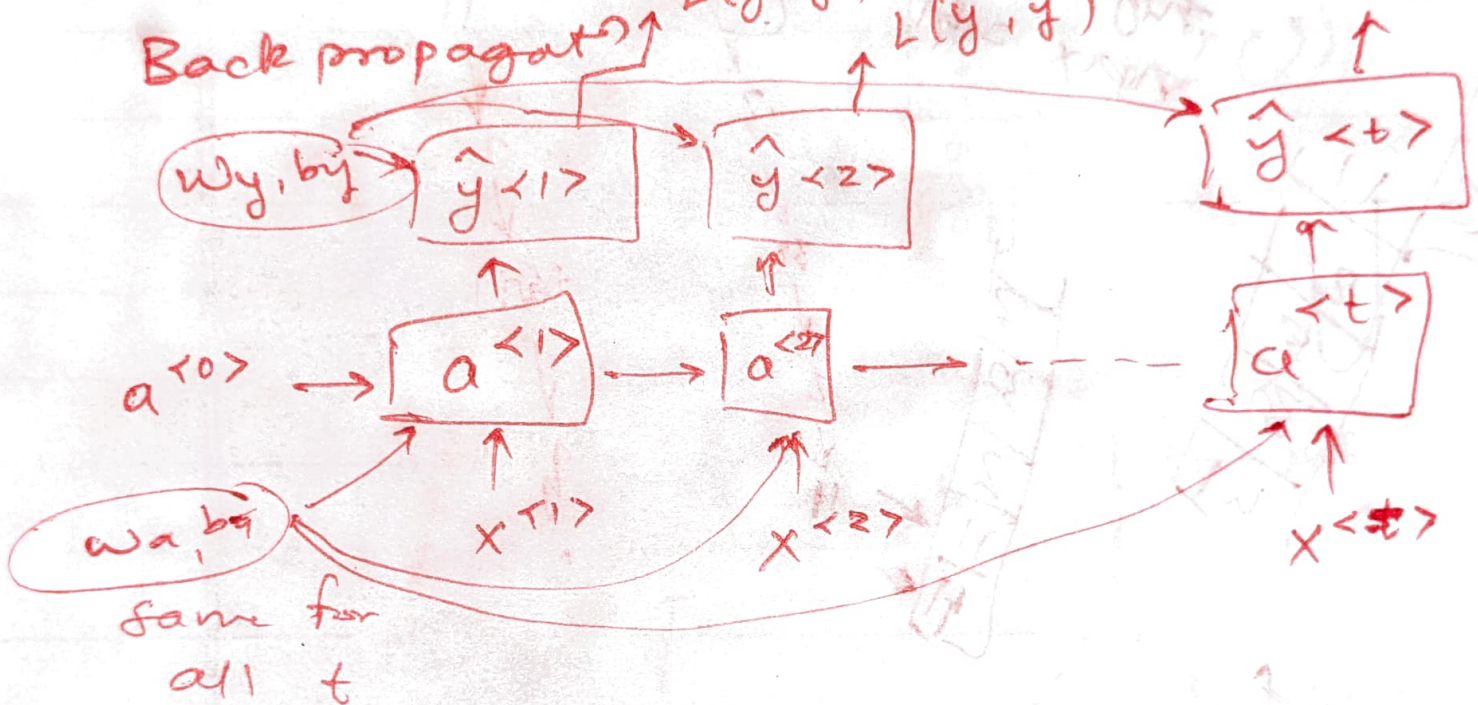
So now we can ~~use~~ use W_a as a single matrix.

∴ we can also write,

$$\hat{y}^{<t>} = g(W_y a^{<t>} + b_y)$$

$L^{<1>}(\hat{y}^{<1>}, y^{<1>})$ $L^{<2>}(\hat{y}^{<2>}, y^{<2>})$ $L^{<t>}(\hat{y}^{<t>}, y^{<t>})$

Back propagation



$$L^{<t>}(\hat{y}^{<t>}, y^{<t>}) =$$

$$-y^{<t>} \log \hat{y}^{<t>} - (1 - y^{<t>})$$

$$\log(1 - \hat{y}^{<t>})$$

for all t ,


$$L(\hat{y}, y) = \sum_{t=1}^T L^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

Types of RNN

- One to One (Typically, not classified as RNN)
- One to many
sound generation
- many to one
sentiment analysis (Movie review to its rating)
- Many to many
 - ① $x^{<t>} \neq y^{<t>}$ are ~~not~~ diff.
eg. translation of diff. lang.
 - ② $x^{<t>}$ & $y^{<t>}$ are same
eg. name entity recognition. Check whether words in a sentence are noun or not.

Language modelling with an RNN.

Language modelling is predicting the next word/output in a sequence of words.

Cats average 15 hours of sleep a _____
predicted word 

Given a dictionary of words, RNN will try to predict the probability of each word.

say dictionary = {a, aaron, ---, zulu},
 \hat{y} i.e. output will consist of probabilities of each word. ~~It~~ And their probabilities sum up to 1.

The probability of each word depends upon the input & the past output. But the predicted output won't be passed as input as it may be incorrect. Hence, during training, the correct output is passed, instead of predicted.

$$P(y^{(t)} | y^{(1)}, y^{(2)}, \dots, y^{(t-1)})$$

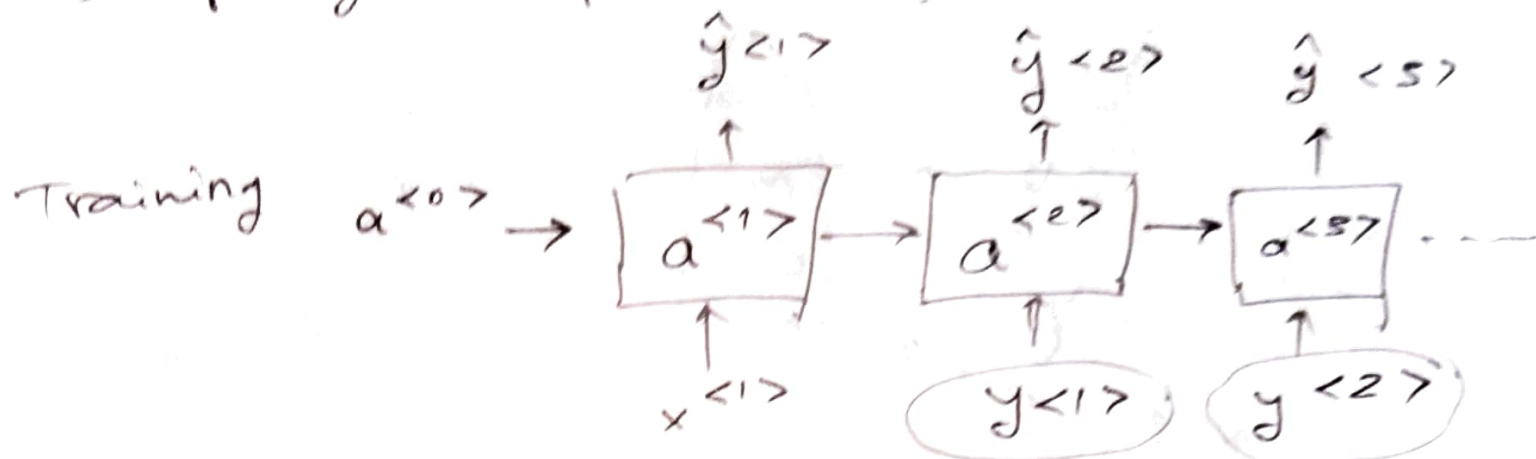
Sampling novel sequences

How it is different from Language modelling!

→ Sampling of a sequences from a trained RNN occurs in the testing phase. A trained model is used to predict the output of each word. We'll just pass an empty input & we'll see, what the model generates from that input.

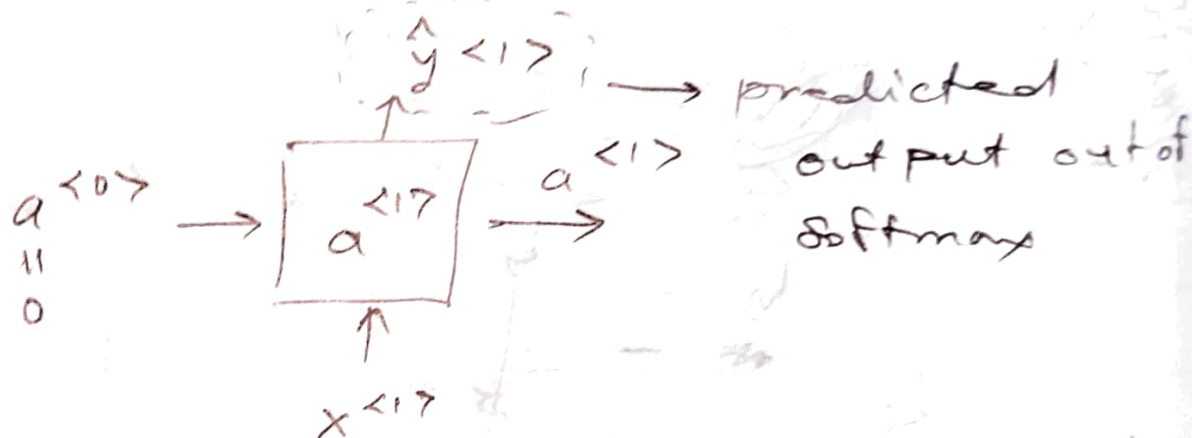
But from a trained model, the output will be always the same. Hence, to avoid that a random word would be chosen from according to their probability. In sampling, the predicted o/p is passed as input to the next layer unlike modelling.

Sampling sequence from trained RNN



Notice, these aren't predicted outputs. They are correct ones (from sample example).

Sampling:



$\hat{y}^{<1>}$ will be consisting of $P(a), P(\text{aaron}),$

.... $P(\text{zulu}), P(\text{<UNK>})$.

Not always, the character / word of highest probability is selected.

np.random.choice is used — ①

while passing $y^{<1>}$ to the next activation function, it is one hot encoded

Meaning all the words except the one which is select by ① is represented by 0 & ① is represented by 1.

When to know to finish complete the sentence?
When $\hat{y}_{<t>}$ will generate $\langle \text{EOS} \rangle$, the end the sentence

Also $\hat{y}_{<t>}$ may generate $\langle \text{UNK} \rangle$ i.e. unknown words. you may discard this result or you may keep it.

2 types of models

① character level language model:
Character by character is generated by each activation function. It can predict unknown words like a person's name which is not included in the vocabulary.

② word level : words are generated from each activation function. Unknown words may not be generated

Vanishing gradients with RNN.
& exploding

if we try to include all the past sequences in the current prediction, even though they are irrelevant, we would get a very large loss at the output. & thus, ~~the~~ during backpropagation, & the gradients would explode. Since, the gradients use chain rule to calculate the gradients in previous layers, the gradients are multiplied. Similarly vanishing of gradient may also occurs, as the gradients may (less than 1) be small, but the gradients might too small during backpropagation making ~~no change~~ almost no changes in the parameters. Hence to avoid this, gradient clipping is done.

The gradients are clipped, i.e., if the ~~gr~~ past gradients are too large, they are normalized to some value, so that all the gradients are in the range of $-max$ to max , given 'max' is the limit.

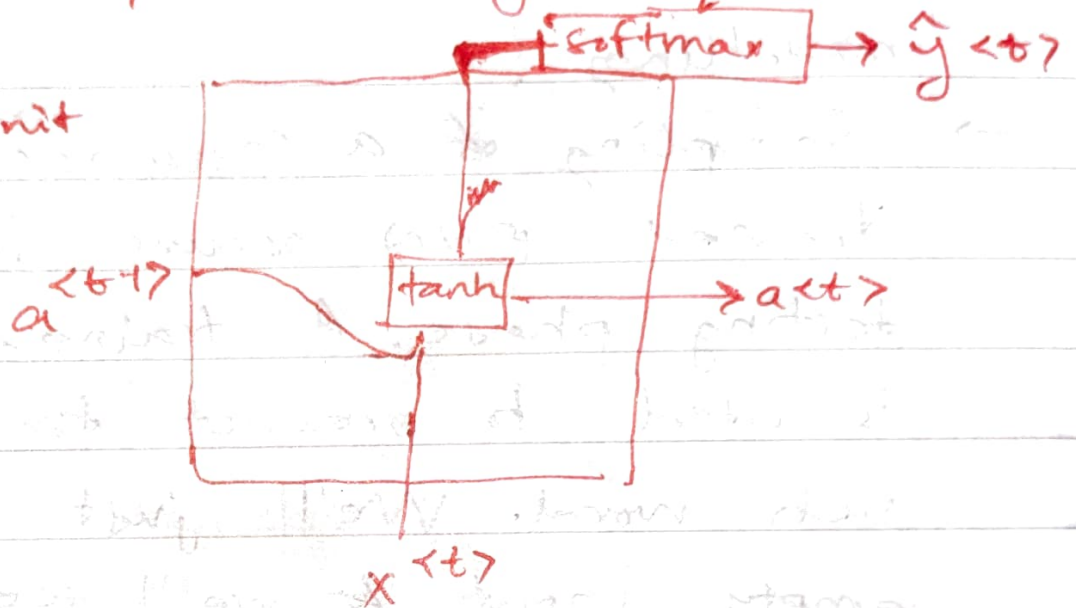
Furthermore, to avoid the ^{very} long range dependencies of words, we'll be using GRU & LSTM.

GRU (Gated Recurrent Unit)

&

We know, $a^{<t>} = g(w_a [a^{<t-1>}, x^{<t>}] + b_a)$

RNN unit



C = memory cell

$C^{<t>} = a^{<t>}$ — for GRU

~~$$\tilde{C}^{<t>} = \tanh(w_c [C^{<t-1>}, x^{<t>}] + b_c)$$~~

$$C^{<t>} = \tanh(w_c [C^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(w_u [C^{<t-1>}, x^{<t>}] + b_u)$$

$u \rightarrow \text{update}$

Γ_u will determine whether the word will be saved for further use or not. It uses a sigmoid function for that

$$\textcircled{1} C^{<t>} = \Gamma_u \times \tilde{C}^{<t>} + (1 - \Gamma_u) \times C^{<t-1>}$$

if Γ_u is 1, ^{update} ~~set~~ gate's value as new value of $C^{<t>}$.

if Γ_u is 0, don't update the gate's value, keep it as it is.

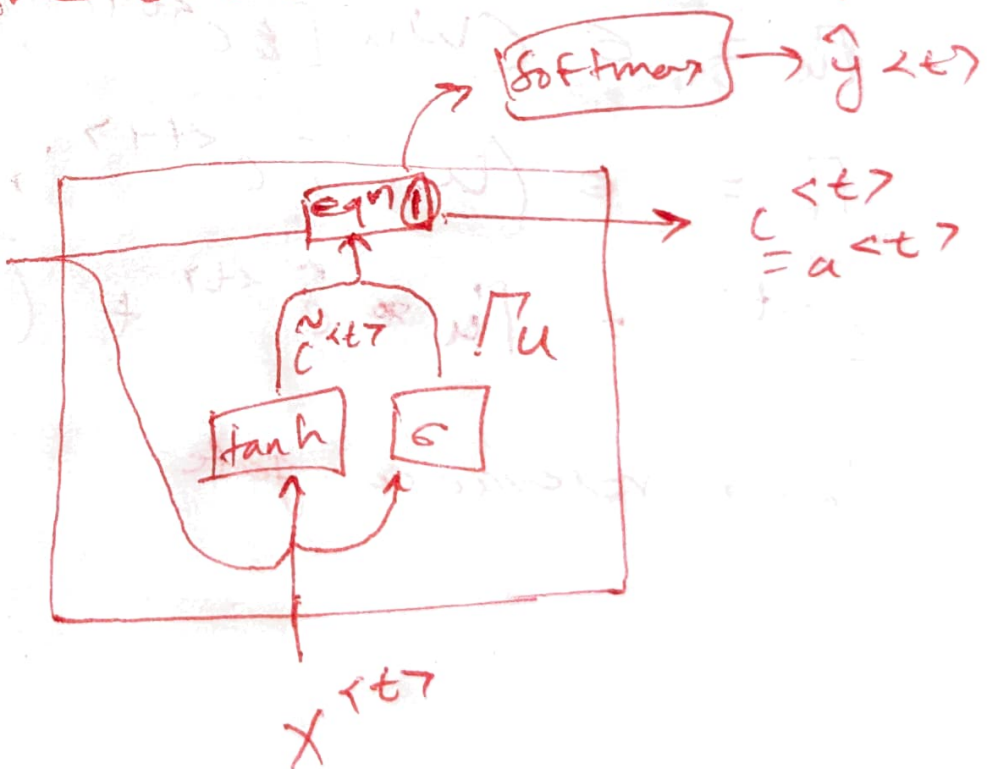
eg.

The cat, which already ~~ate~~ ate \longrightarrow was full.

\longrightarrow for cat, ~~to~~ update gate's value to 1, & for all the other words, keep it 0.

GRU

$$C^{<t-1>} = a^{<t-1>}$$



Γ_u doesn't suffer much from vanishing gradient.

$c^{<t>}$, $\tilde{c}^{<t>}$, Γ_u will have same dimensions and would depend on the dimension of activation functions.

$a^{<t>}$

In $c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$

element wise multiplication

* Full GRU

$$\tilde{c}^{<t>} = \tanh(W_c [\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u [\tilde{c}^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r [c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$\Gamma_r \rightarrow$ relevance gate

LSTM

$$\tilde{c}^{<t>} = \tanh [w_c [a^{<t-1>}, x^{<t>}] + b_c)$$

You may include Γ_r in above eqn as well, but more commonly its not part of the eqn

update $\Gamma_u = \sigma (w_u [a^{<t-1>}, x^{<t>}] + b_u)$

forget $\Gamma_f = \sigma (w_b [a^{<t-1>}, x^{<t>}] + b_f)$

output $\Gamma_o = \sigma (w_o [a^{<t-1>}, x^{<t>}] + b_o)$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

~~$$a^{<t>} = \Gamma_o * c^{<t>}$$~~

$$a^{<t>} = \Gamma_o * \tanh c^{<t>}$$

Bidirectional RNN

Traditional RNN & LSTM only have ~~relevance~~ ^{relevance} info from the past & no ~~relevance~~ ^{relevance} info from the future

- ① He said, "Teddy bears are on sale!"
- ② He said, "Teddy Roosevelt was a great President".

Here, each "Teddy" has different meaning. And we cannot decide it only from the information from the past.

BRNN requires entire sequence to make prediction. Not helpful for real time application. ~~where~~ say for speech recognition, we would require to complete our sentence for BRNN to make ~~sequence~~ ^{sequence} of words product.